

Information Theoretic Approaches to Rapid Discovery of Relationships in Large Climate Data Sets

Kevin H. Knuth¹, William B. Rossow²

1. NASA Ames Research Center, Moffett Field CA, 94035-1000

2. NASA Goddard Institute for Space Studies, New York NY, 10025

Mutual information as the asymptotic Bayesian measure of independence is an excellent starting point for investigating the existence of possible relationships among climate-relevant variables in large data sets. As mutual information is a nonlinear function of its arguments, it is not beholden to the assumption of a linear relationship between the variables in question and can reveal features missed in linear correlation analyses. However, as mutual information is symmetric in its arguments, it only has the ability to reveal the probability that two variables are related. It provides no information as to how they are related; specifically, causal interactions or a relation based on a common cause cannot be detected. For this reason we also investigate the utility of a related quantity called the transfer entropy. The transfer entropy can be written as a difference between mutual informations and has the capability to reveal whether and how the variables are causally related. The application of these information theoretic measures is tested on some familiar examples using data from the International Satellite Cloud Climatology Project (ISCCP) to identify relations between global cloud cover and other variables, including equatorial pacific sea surface temperature (SST), over seasonal and El Nino Southern Oscillation (ENSO) cycles.

Effective Approaches to Rapid Discovery of Relationships in Large Climate Data Sets

Kenneth K. Kuhl, William B. Rosow

and the fact that the fundamental frequency of vibration varies from 100 to 150 Hz. According to the ISO 2631-1997 standard, the maximum allowable value is 0.25 ms⁻² r.m.s.

B61A-0703

15122416N

Mutual information as the asymptotic Bayesian measure of interdependence is an excellent starting point for investigating the existence of possible relationships among climate-relevant variables in large data sets. As mutual information is a nonlinear function of its arguments, it is not beholden to the assumption of a linear relationship between the variables in question and can reveal features missed in linear correlation analyses. However, as mutual information is asymptotic in its arguments, it only has the ability to reveal the probability that two variables are related. It provides no information as to how they are related, specifically, causal interpretations or a relation based on a common cause cannot be detected. For this reason we also investigate the utility of a related quantity called the transfer entropy. The transfer entropy can be written as a difference between mutual informations and has the capability to reveal whether and how the variables are causally related. The application of these information theoretic measures is illustrated on some familiar examples using data from the International Satellite Cloud Climate Project (ISCCP) to identify relationships between global cloud cover and other variables, including equatorial sea surface temperature (SST), over seasonal and El Niño Southern Oscillation (ENSO) cycles.

WOMAN'S INSTITUTION

We can characterize the behavior of a system X by looking at the probability distribution over the set of states the system visits as it evolves in time. If a state is visited rarely, we might be surprised to find the system there. We can express the expectation (or lack of) to find the system in state x as

$$h(x) = \log \frac{1}{p(x)}$$

If we average this over all states, this gives us a measure of our expectation (or our uncertainty). This quantity is called the Shannon entropy [4]

$$H(X) = - \sum_{x \in X} p(x) \log \frac{1}{p(x)}$$

If the system states can be described with multiple parameters, the entropy can still be computed by averaging over all possible states (here is shown for states described by X and Y)

$$(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{1}{p(x, y)}$$

Now, if we consider two subsystems X and Y , which together make up a larger system, we can compute what is called the **Mutual Information (MI)** by

$$MI(X, Y) = (H(X) + H(Y)) - H(X, Y)$$

Notice that this describes the difference between the uncertainty when the two are treated separately and when they are treated jointly. If these two sub systems are independent of one another, the MI will be zero. However, if there is any interaction between these subsystems, the MI will be positive. This can perhaps be seen more clearly by writing it as

$$M(X, Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

X and Y are independent then the probability in the numerator factors and the MI is zero.

[illegible]

However, as the MI is symmetric with respect to interchange of X and Y it not distinguish the direction of any causal influence as well as the act of any common influence on both subsystems from a third system. A similar quantity called the **Transfer Entropy** can be written as a difference of MIs

$$T(X_{i+1} | X_{i+1}^{(0)}, X_{i+1}^{(n)}) = \lambda \Pi(X_{i+1}, X_{i+1}^{(n)} \otimes Y_{i+1}^{(n)}) - \lambda \Pi(X_{i+1}, X_{i+1}^{(0)})$$

where $X_i^{(0)}$ is the joint system defined by the set of states from X_{i+1} to X_i , $X_i^{(0)} \otimes Y_i$ is the system represented by the k previous states from X_i with the k previous states from Y . Thus it measures the change in knowledge obtained from incorporating information about the subsystem for $k \ll \infty$. This can be nicely written in terms of Shannon entropy (though not as clearly interpretable)

$$X_{i,\sigma} | X_i^{(1)}, Y_i^{(n)} = -H(X_i) + H(X_i, Y_i) + H(X_i, Y_j) - H(X_i, Y_{i+1}, Y_j)$$

Reasonability and Religious Conviction

We used data (climate summary product C2) from the International Satellite Cloud Climatology Project (ISCCP) to investigate the ability of these information-theoretic techniques to efficiently and accurately discover relationships between seasonal change and global cloud cover.

The data consisted of monthly averages of percent cloud cover resulting from a time-series of 198 months of 6596 equal-area pixels each with side length of 280 km. The analysis was performed pixel-wise so that for each pixel: X = cloud cover percentages and Y = month of the year (seasonal date). The MI was computed for each pixel independently and is color-coded on the map below.



This method finds the Inter-Tropical Convection Zones, The Monsoon regions, the Sea Ice off Antarctica, and cloud cover in the North Atlantic and Pacific. This figure can be directly compared to the PCA analysis performed by Rossow et al. [5].

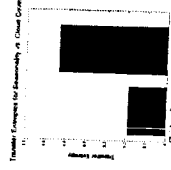
ocusing on the high-MI region in the Katanga Plateau of southern ongo we see that the joint probability density is not factorable as indicated by the MI. The summer months are sunny and the winter onths are cloudy.



In contrast the cloud cover over Paris France has a low MI with respect to the seasonality. This is reflected by the highly factorable joint probability density to the right. Note that there is little dependence of cloud cover on seasonality.



he transfer entropies were found to be difficult to estimate with precision. Using a Gaussian kernel density estimator, we found that the seasonality over the cloud cover, which is reassuring as we know that the seasonality is an autonomous variable.



NSO and Clonidine

The next example looks at the relationship between global cloud cover and the Cold Tongue Index (CTI) [6], which describes the sea surface temperature anomalies in the eastern equatorial Pacific Ocean (6N-6S, 180-90W) and is indicative of ENSO variability [7]. Data is from [8].

The cloud cover affected by the SST variations lies mainly in the equatorial Pacific, along with an area in Indonesia. The highlighted areas in the Indian longitudes are artifacts of satellite coverage.



the probability density below shows how the two variables are co-dependent, whereas the transfer entropy indicates the direction of interaction.



Entropy estimates can be improved using Markov chain Monte Carlo. This approach will result in estimates with error bars. In addition, we will be able to directly test for independence of climate parameters.

15

- Wolf D R. 1994. Mutual information as a Bayesian measure of independence. *Le UPR Vol 9*, 49-60.
- Schreiber T. 2000. Measuring information transfer. *Phys Rev Lett* 85:461.
- Kaiser A, Schreiber T. 2002. Information transfer in continuous processes. *Physica D* 166:13-32.
- Shannon CE. Weaver W. 1949. *The Mathematical Theory of Information*. University of Illinois Press, Urbana IL.
- Rosenau WB, Walter AM, Grider LC. Comparison of ESCTP and other cloud microphysics. *J Climate* 13:2411-2441.
- Dwyer C, Wallace JM. 1990. Large-scale atmospheric circulation features of warm and cold episodes in the tropical Pacific. *J Climate* 3:1234-1281.
- Dwyer C, Wallace JM. 1997. El Niño events and their relation to the Southern Oscillation. 1925-1986. *J Geophys Res* 92:14189-14196.
- <http://www.washington.edu/education/understandments/1809/05.html>

Funded by
Intelligent Data Understanding Project / Intelligent Systems Program
ASA Aerospace Technology Enterprise